

# PATENT ABSTRACTS OF JAPAN

(11)Publication number: 10-293860

(43)Date of publication of application: 04.11.1998

(51)Int. Cl. G06T 13/00  
G10L 3/00

(21)Application number: 09-216377 (71)Applicant: NIPPON TELEGR & TELEPH CORP <NTT>

(22)Date of filing: 11.08.1997 (72)Inventor: MATSUURA MICHIAKI  
KURA TSUNEKO  
OSHIMA TAKASHI  
WATANABE NOBUYUKI  
KANAYAMA HIDEAKI

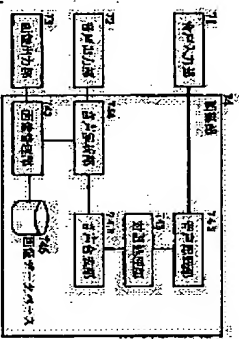
(30)Priority

Priority number: 09 38760 Priority date: 24.02.1997 Priority country: JP

(54) PERSON IMAGE DISPLAY METHOD AND DEVICE USING VOICE DRIVE

(57)Abstract:

PROBLEM TO BE SOLVED: To artificially attain an almost natural conversation state with no use of any special device, etc., by changing at least one of shakes of a mouth, head, etc., of a person face image based on the loudness of human voices which are obtained in sequence and at each fixed time interval. SOLUTION: A voice analysis part 744 picks up the changes of voices received from a voice synthesizing part 743 at each prescribed time interval and calculates the mean value of loudness of voices, etc. Then the time series information on each of mean value of voice loudness divided at each time interval of voices is outputted to an image management part 745. The part 745 successively selects the prescribed proper images out of an image data base in response to each information and outputs these images to an image output part 73 to successively display them. These images primarily show an answering person to the voices of a speaker, and the mouths, attitudes, etc., of one or more persons are operated in accordance with the voices produced from a voice output part 72.



(19) 日本国特許庁 (JP)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平 10-293860

(43) 公開日 平成 10 年 (1998) 11 月 4 日

(51) Int. Cl. G 06 T 13/00  
G 10 L 3/00

F I G 06 F 15/62 3 4 0 D  
G 10 L 3/00 S

審査請求 未請求 請求項の数 11 OL

(全 14 頁)

(21) 出願番号 特願平 9-216377

(22) 出願日 平成 9 年 (1997) 8 月 11 日

(31) 優先権主張番号 特願平 9-38760

(32) 優先日 平 9 (1997) 2 月 24 日

(33) 優先権主張国 日本 (JP)

(71) 出願人 000004226  
日本電信電話株式会社  
東京都新宿区西新宿三丁目 19 番 2 号

(72) 発明者 松浦 道明  
東京都新宿区西新宿三丁目 19 番 2 号 日本  
電信電話株式会社内

(72) 発明者 倉 恒子  
東京都新宿区西新宿三丁目 19 番 2 号 日本  
電信電話株式会社内

(72) 発明者 大島 孝  
東京都新宿区西新宿三丁目 19 番 2 号 日本  
電信電話株式会社内

(74) 代理人 井理士 秋田 収彦  
最良頁に続く

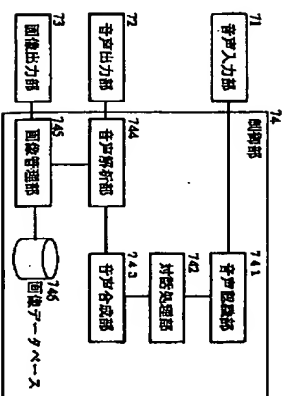
(54) 【発明の名称】 音声駆動を用いた人物画像表示方法およびその装置

(57) 【要約】

【課題】 専門の技術あるいは装置を必要となくともほぼ自然的に対話をしている状態を駆動的に実現できる。

【解決手段】 コンピュータに表示されている人物の画像の少なくとも顔の部分にその人物の発する音声によって変化する音声駆動を用いた人物画像表示方法において、前記人物の発する音声の一定時間間隔ごとに順次得られる各音声の大きさに基づいて、該人物の顔画像における口、頭の揺れあるいはうなずき、目のまばたきの少なくとも一つを変化させる。

図 1





求める。

【0024】さらに、音声の大きさの各平均値 $S_a$ に基  
づいて、図3に示すように、前記各画像の遷移先を決定  
する。

【0025】図3においては、4種類の口の大きさを持  
った画像を明示しており、たとえば音声入力がない場合  
にはすべて口を閉じた画像に遷移する。音声の大きさの  
最大値を4つに分け、その各々の区間に口の大きさの異  
なる画像を配置する。 $S_a$ の値が変わることにより、口  
の大きさも変化する。このとき、口の大きさは急激に変  
わる場合もある。

【0026】そして、拡張アルゴリズムとして、前記基\*

$$P_{B=0} + P_{B=1} + P_{B=2} = 1 \quad \dots\dots\dots (1)$$

$P_{B=0}$ : 基準となる顔から基準となる顔へ遷移する確率

$P_{B=1}$ : 基準となる顔から顔を揺らした顔へ遷移する確率

$P_{B=2}$ : 基準となる顔から目を閉じた顔へ遷移する確率

【0030】なお、この場合、上式(1)の各確率はそ  
の等式が満足できる限りどのような割合にしてもよいこ  
とはいまでもない。

【0031】ここで、前記基本アルゴリズムについてさ  
らに詳述すると、上述したように、ある一定時間間隔 $T$   
ごとに音声の平均値 $S_a$ を求め、さらに、図5に示す  
ように、ある時間 $T_n$ と $T_{n+1}$  ( $=T_n+T$ ) との平均値  
の差分を $D (=S_{a,n+1}-S_{a,n})$  とおく。

【0032】そして、この差分 $D$ の大きさにより、図6  
に示すように、画像の遷移先を決定する。同図では、差  
分 $D$ がたとえあればあるしきい値 $\epsilon$ 内に収まる場合には自分  
自身に遷移するようにしている。差分 $D$ がしきい値 $\epsilon$ を  
越え、かつその値が正の場合には口をさらに大きく開け\*

$$P_{B=0} + P_{B=1} + P_{B=2} = 1 \quad \dots\dots\dots (2)$$

$P_{B=0}$ : 基準の顔で口を開いた顔へ遷移する確率

$P_{B=1}$ : 顔を揺らした口を開いた顔へ遷移する確率

$P_{B=2}$ : 目を閉じて口を開いた顔へ遷移する確率

【0036】そして、差分 $D$ が負の場合には第 $(n-1)$  状態への遷移は次式(3) に示す等式を満足して行  
われるようになっている。

$$P_{B=0} + P_{B=1} + P_{B=2} = 1 \quad \dots\dots\dots (3)$$

$P_{B=0}$ : 基準の顔で口を開いた顔へ遷移する確率

$P_{B=1}$ : 顔を揺らした口を開いた顔へ遷移する確率

$P_{B=2}$ : 目を閉じて口を開いた顔へ遷移する確率

【0038】差分 $D$ がしきい値 $\epsilon$ 内にある場合には自分  
自身へ遷移し、また、音声の大きさが0の場合には第1  
状態に遷移することはいまでもない。

【0039】このようにすることによって、しきい値 $\epsilon$  50

を致す実験を行った結果を以下に説明する。

【0041】この実験における基本アルゴリズムとし  
て、図8 (a) に示すように、声の大きさの変化にとも  
なって口の大きさが4段階に変化する顔画像を用意し  
た。

【0042】そして、拡張アルゴリズムとして、図8  
(b) に示すように、上述した3パターンに、目を半分  
閉じたパターンを追加し、遷移先は一つ前の状態には戻  
らない。数定とした。また、音声を受信している際には、  
瞬きをしたり、うなずく等の動作をランダムに表示する  
ようにした。

表示アルゴリズム評価結果  
Table 2 The Results of Algorithms

カテゴリ	単純尺度値	基本アルゴリズム	拡張アルゴリズム
非常によい	4	1	1
よい	3	1	3
普通	2	2	4
悪い	1	5	2
非常に悪い	0	1	0
利便の可視尺度値		1.60	2.30

【0046】この表1から明らかなように、拡張アル  
ゴリズムによる評価 (尺度値=2、3) が基本アル  
ゴリズムの評価 (尺度値=1、6) より高いことが判る。

【0047】この理由は、基本アルゴリズムでは、顔  
他の部分を動かさなかったため、口の動きだけが目だっ  
てしまい、時間を進めるに従ってその不自然さが増大する  
状態に陥り易いからである。拡張アルゴリズムではこの  
点の不都合を充分に解決したものとされている。

【0048】また、表示されている人物が複数人おり、  
そのうちの一人の話を聞いている人物の場合において、  
その口を閉じた状態の画像で、目を閉じた顔と目を開い  
た顔との表示をランダムに切り替えることによって、そ  
の人物の目の瞬きを表現できるようにしている。

【0049】このように構成した人物画像表示方法およ  
び装置の利用法としては、たとえば商品紹介、タレント  
イド、あるいは交通情報表示等の情報提供が挙げられ  
る。

【0050】たとえば交通情報の場合は、具体的に、交  
通手段 (鉄道/バス/道路/飛行機...)、混み具合 (所  
要時間/混雑距離/混雑原因...)、その他の情報 (天気  
予報/工事情報/事故...) 等である。

【0051】利用者は、たとえば、提示されているメニ  
ューから必要な項目を音声で入力することになる。入力  
された音声は音声認識部で認識され、認識された言葉に  
対応する返答を合成/解析し、その結果に基づき画面に  
表示する顔を動かす。

\* 【0043】そして、被験者10人に対して、それぞれ  
基本アルゴリズムに基づく画像と拡張アルゴリズムに基  
づく画像とを提示し、それにより得られた結果を評定尺  
度法 (大甲、中山、堀田、" 画質と音質の評価技術" 昭  
和堂: 1991) に基づいて評価してみた。

【0044】すなわち、カテゴリ-数を5とし、単純尺  
度値により各々の判断の尺度値を求め、その結果は表1  
に示されるようになった。

【0045】

【表1】

【0059】まず、対話というものは、一方的に言葉を

相手に投げけるものでなく、相手の反応を見ながら行うものである。

【0060】対話がスムーズになされている状態を表現するためには、話をしている側の顔の部分にそれなりの動きをもたせることはもちろんのこと、聞いている側の顔の部分もそれなりに表現させる必要がある。

【0061】そこで、対話の際の話をしている側および聞いている側のそれぞれの顔の動きを観察すると次のような動作傾向があることが判る。

- ・顔が動くことが多い。
- ・瞬きの回数が一回が多い。
- ・瞬きから瞬きの間隔が比較的長い。
- 【0063】(2)話を聞いている側の動作
  - ・瞬きは複数回連続して行われることが多い。
  - ・相槌をうつため、縦方向に顔が動く。
  - ・疑問を感じた際に、首を傾げる。
  - ・納得できない時に、首を振る。
- 【0064】図9および図10に示す実施例では、対話者のこのような動作傾向を出力制御に反映させていることにある。

【0065】すなわち、コンピュータに表示されている人物が、話をしているのか(有音)、あるいは話を聞いているのか(無音)の判定を行うとともに、それぞれの場合における顔の部分の動作を、口の大きさ、目の瞬きに着目して制御するようにしている。

【0066】図9、図10では、図面の横方向に時間経過を示し、上段から順次、①有音/無音の判定タイミング、②有音/無音出力、③口画像出力、④目画像出力、⑤顔画像出力を示している。そして、図9は、主として有音判定時を示し、図10は、主として無音判定時を示している。なお、図9、図10は、それらが合わさって一つの図を構成するようにしている。

【0067】ここで、①有音/無音の判定タイミングでは、有音/無音の判定が一定周期 $T_1$ でなされている。

【0068】②有音/無音出力では、有音判定時に、入力音が無音レベルから有音レベルに変化したことにより、無音から有音に変化することを示し、また無音判定時に、入力音が無音レベルから無音レベルに変化したことにより、有音から無音に変化することを示している。

【0069】有音/無音の判定は、たとえば判定間隔 $T_1$ の中に存在する標準化音声の平均電力を計算し、その電力が予め定めたしきい値を超えるかどうかで行うようにしている。

【0070】③口画像出力では、有音判定時に、有音が開始してから一定周期( $T_m$ )で口の大きさが変化することを示し、また無音判定時に、直ちに口が閉じることを示している。

【0071】④目画像出力では、有音判定時に、有音が

(6)

10

特開平10-293860

開始してからランダムに瞬きが発生し、その後は一定の周期( $T_m$ )で瞬きが繰り返すことを示し、また無音判定時に、無音が開始してからランダムに瞬きが発生し、その後は一定の周期( $T_m$ )で瞬きが繰り返すことを示している。

【0072】この場合、無音では、瞬きの回数が2回連続して行われることを示している。

【0073】⑤顔画像出力、口画像と目画像を合成した結果を示している。

【0074】なお、上述したタイムチャートでは、有音から無音へ移行した段階で、直ちに口を閉じるようにして示したが、このようにすることに限定されることはなく、検出時から時間をかけて漸次に口を閉じる制御を行うようにしてもよいことはいずれまでもない。

【0075】また、口画像出力において、口の閉閉の繰返し周期 $T_m$ を一定として説明したものであるが、このようにすることに限定されることはなく、ランダムの時間幅で口の閉閉を繰り返すようにしてもよいことはいずれまでもない。

【0076】さらに、目画像出力において、目の瞬きの周期 $T_m$ 、 $T_m'$ を一定として説明したものであるが、このようにすることに限定されることはなく、ランダムの周期で目の瞬きを繰り返すようにしてもよいことはいずれまでもない。

【0077】図11および図12は、図9および図10に示したタイムチャートを上述した人物画像表示装置に実行するためのフロー図である。なお、図11、図12は、それらが合わさって一つの図を構成するようにしている。また、このフローの実行に先立ち、図13に示す口変数デューツル(1)、目変数デューツル(2)を作成しておき、画像データベース746に格納しておく。

【0078】以下、ステップ順に説明する。

【0079】ステップ1(S1)は、まず、初期段階として、口の変数 $mouth=0$ 、目の状態変数 $eye=0$ としておく。

【0080】ステップ2(S2)は、 $mouth=0$ 、 $eye=0$ にそれぞれ対応する口画像および目画像を出力する。この場合、それぞれの画像は閉じた口および閉じた目となっている。

【0081】ステップ3(S3)は、入力音を取り入れ、それが有音であるか無音であるかを判定する。有音である場合はステップ4へいき、また、無音である場合はステップ14へいき。

【0082】ステップ4(S4)は、ステップ3で判定した有音が2度目以上の有音が初めて有音であるかを判定する。初めての有音である場合はステップ5へいき、また、2度目以上の有音である場合はステップ6へいき。

【0083】ステップ5(S5)は、変数 $i$ に乱数装(0、 $T_1$ 、 $2T_1$ 、……、 $T_m$ )から選

11

(7)

特開平10-293860

ばれた数値を代入する。瞬きの開始時間をランダム化させるためである。

【0084】ステップ6(S6)は、変数 $i$ を $T_m$ で割り算を実行し割り切れるか否かを判定する。瞬きを $T_m$ 毎に行わせるためである。割り切れる場合はステップ7へいき、割り切れない場合はステップ8へいき。

【0085】ステップ7(S7)は、 $eye=1$ に相当する目の画像を出力する。そして、変数 $i$ を $T_1$ だけ増加させる。

【0086】ステップ8(S8)は、 $eye=0$ に相当する目の画像を出力する。そして、変数 $i$ を $T_1$ だけ増加させる。

【0087】ステップ9(S9)は、口の大きさが増加中か否かを判定する。増加中の場合はステップ10へいき、そうでない場合はステップ12へいき。

【0088】ステップ10(S10)は、口の大きさが最大値に達したか否かを判定する。達した場合にはステップ13へいき、達しない場合にはステップ11へいき。

【0089】ステップ11(S11)は、口の変数に1を加算し、その加算値に対応する口画像を出力する。その後、ステップ3へいき。ステップ10からステップ11までの動作は、口の大きさが増大しているのであれば、その最大値まで口の大きさを増大させるためである。

【0090】ステップ12(S12)は、口の大きさが最小値に達したか否かを判定する。達した場合にはステップ11へいき、達しない場合にはステップ13へいき。

【0091】ステップ13(S13)は、口の変数に1を減算し、その減算値に対応する口画像を出力する。その後、ステップ3へいき。ステップ12からステップ13までの動作は、口の大きさが減少しているのであれば、その最小値まで口の大きさを減少させるためである。

【0092】ステップ14(S14)は、2度目以上の無音が初めての無音か否かを判定する。初めての無音の場合はステップ15へいき、2度目以上の無音の場合はステップ16へいき。

【0093】ステップ15(S15)は、変数 $i$ に乱数装( $T_1$ 、 $2T_1$ 、……、 $T_m'$ )から選ばれた数値を代入する。瞬きの開始時間をランダム化させるためである。

【0094】ステップ16(S16)は、変数 $i$ を $T_m'$ で割り算を実行し割り切れるか否かを判定する。瞬きを $T_m'$ 毎に行わせるためである。割り切れる場合はステップ17へいき、割り切れない場合はステップ18へいき。

12

特開平10-293860

【0095】ステップ17(S17)は、 $eye=1$ 、0、1となる目の画像を順次出力する。それらの出力は $T_1$ 間隔となり、その後において、 $i$ に3 $T_1$ を加算する。

【0096】ステップ18(S18)は、 $eye=0$ の目の画像を出力する。その後、 $i$ に $T_1$ を加算する。

【0097】ステップ19(S19)は、無音期間中は $mouth=0$ として口を閉じておく。そして、ステップ3へいき。

【0098】実施例3. 図14は、上述した音声駆動を用いた人物画像表示装置をネットワークを介して他の音声駆動を用いた人物画像表示装置とインテラティブに接続を行う場合の実施例を示したブロック図である。

【0099】同図において、図1と同符号のものは同一の機能を有する。

【0100】図1と異なる構成は、一方の人物画像表示装置側の対話者が発する音は音声入力部71を介してデータ送信部7431に送られ、さらにネットワーク75を介して他方の人物画像表示装置側のデータ受信部7432に送られるようになっている。

【0101】そして、他方の人物画像表示装置の制御部では、データ受信部7432によって受けた音声信号から、図1の場合と同じように、その言葉の意味を判断し、対話処理部742によって対話者の意図を成り取り、かつ、それに対する返答を判断するようになっている。

【0102】音声解析部744では、データ受信部7432からの音に対してたとえば $T_1$ の時間間隔ごとに該音声を変化を取りだし、その強さの平均値の算出、あるいはその他の必要な演算処理を行うようになっている。

【0103】そして、この音声解析部744で得られた音声の時間毎に区別された各平均値等の時系列的な情報は画像管理部745に出力され、この画像管理部745では、前記それぞれの情報に応じて画像データベース746から予め定められた適当な画像が順次選択され、この画像は画像出力部73に出力されて順次表示されるようになっている。

【0104】この画像は、前記対話者の音に対して返答をする物の画像(一人の顔もあるし、複数の顔もある)が主となり、その人物の口あるいはその手(同じタイミングで)動くようになっている。

【0105】この場合において、ネットワークを介して接続されるそれぞれの人物画像表示装置における音声解析部744、画像管理部745、および画像データベース746は、実施例1にて示した基本アルゴリズムおよびそれをさらに拡張させた拡張アルゴリズムに基づいて機能することはいずれまでもない。

【0106】また、このような構成は、たとえばインターネットフオンのように、パソコンで通信しながらお互いの音声を聞いて会話するようにすることもできる。この場合、音声だけでなく、画像データベース746に登録してある低画像（たとえば唇をする相手の低画像）を自分のパソコンの画面上に表示する。データ受信部7432から送られてきた音声の強弱などの情報を音声解析部744で求め、その結果を受けて画像データベース746から適当な画像が順次選択され、画像出力部73に出力される。

【0107】以上、上述した各実施例によれば、人物の表情において、たとえばその口が閉じるようになり、または唇の大きさが等分された数に近づいてそれらの間で口の開きの大きさを変化させることができることから、ほとんどの対話者が不自然さを感じることなく観察することができるようになる。

【0108】したがって、特に、複雑な手法を用いる必要がないことから、専門の技術あるいは装置を必要とすることなくほぼ自然的に対話をしている状態を模似的に実現できるようにする。

【0109】なお、上述した各実施例では、いずれも基本アルゴリズムおよび拡張アルゴリズムにしたがって動作するように説明したものであるが、これに限定されることはなく、それらの一方にしかたがって動作できるようにしてもよいことはいふまでもない。

【0110】

【発明の効果】以上説明したことから明かなように、本発明による音声駆動を用いた人物画像表示方法および装置によれば、専門の技術あるいは装置を必要となく、ほぼ自然的に対話をしている状態を模似的に実現できることになる。

【図面の簡単な説明】

【図1】本発明による音声駆動を用いた人物画像表示装置の一実施例を示すブロック図である。

【図2】本発明の基本アルゴリズムの一実施例を示す説明図である。

【図3】本発明の基本アルゴリズムの一実施例を示す説明図である。

【図4】本発明の基本アルゴリズムの一実施例を示す説明図である。

【図5】本発明の拡張アルゴリズムの一実施例を示す説明図である。

【図6】本発明の拡張アルゴリズムの一実施例を示す説明図である。

【図7】本発明の拡張アルゴリズムの一実施例を示す説明図である。

【図8】本発明の拡張アルゴリズムの一実施例を示す説明図である。

【図9】本発明の拡張アルゴリズムのさらに拡張された実施例を示すタイムチャートで、図10とで一つの図を構成している。

【図10】本発明の拡張アルゴリズムのさらに拡張された実施例を示すタイムチャートで、図9とで一つの図を構成している。

【図11】図9および図10に示すタイムチャートをフロー図で、図12とで一つの図を構成している。

【図12】図9および図10に示すタイムチャートをフロー図で、図11とで一つの図を構成している。

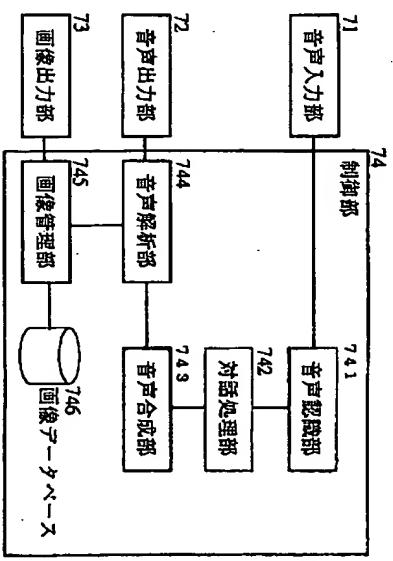
【図13】図11および図12に示すフロー図の動作に先立って予め装置に格納させておくデータを示す説明図である。

【図14】本発明による音声駆動を用いた人物画像表示装置の他の実施例を示すブロック図である。

【符号の説明】

71……音声入力部、72……音声出力部、73……画像出力部、74……制御部、741……音声認識部、742……対話処理部、743……音声合成部、744……音声解析部、745……画像管理部、746……画像データベース。

【図1】



【図13】

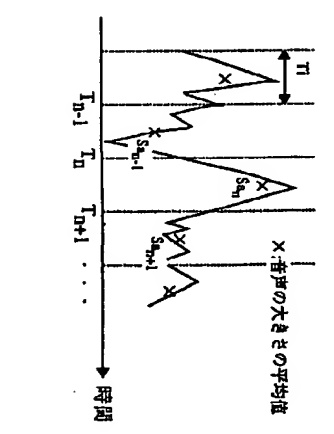
口変換の値と口開度

口変換 (mm)	口開度
0	—
1	○
2	○
3	○
...	...
N	○

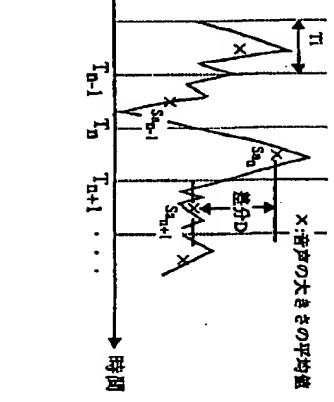
口の位置と口の開口度

口の位置 (mm)	口の開口度
0	—
1	○

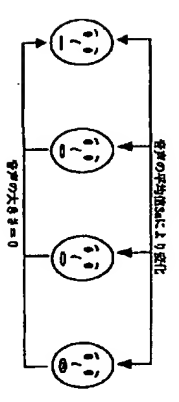
【図2】



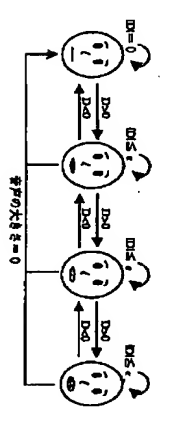
【図5】



【図3】



【図6】



【図4】

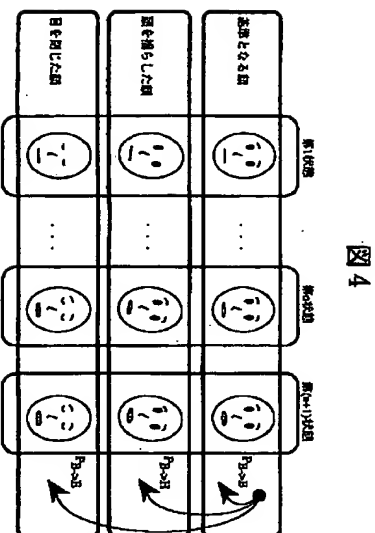


図4

【図8】

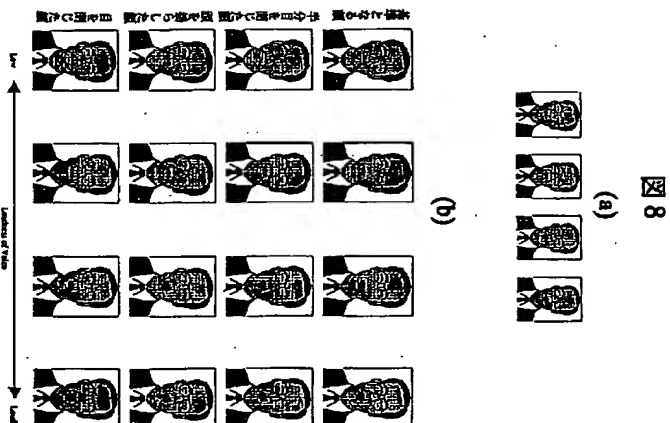
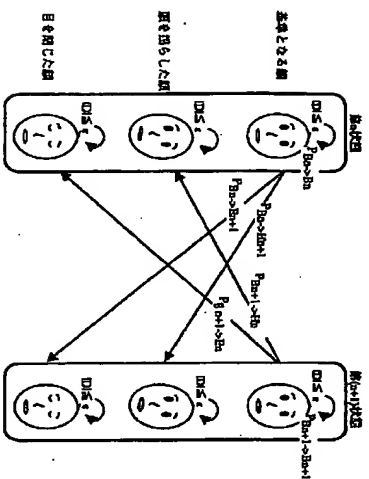


図8

【図7】

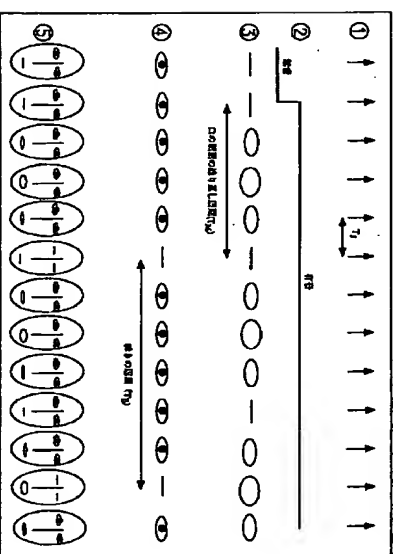
図7



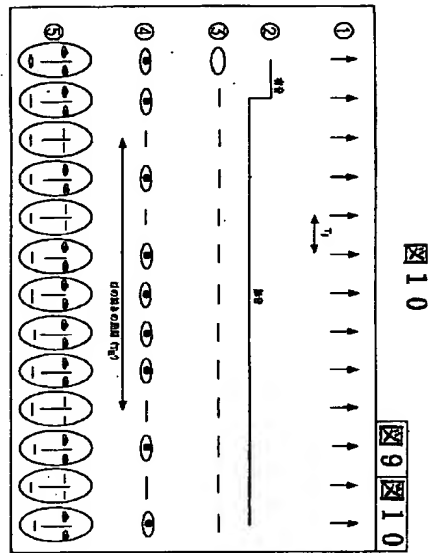
【図9】

図9

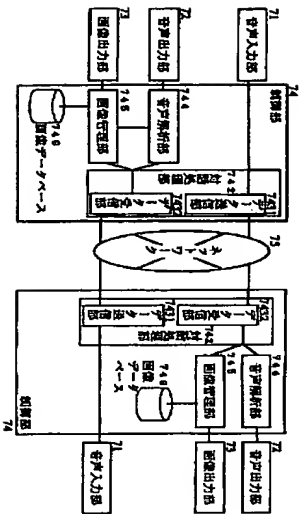
図9図10



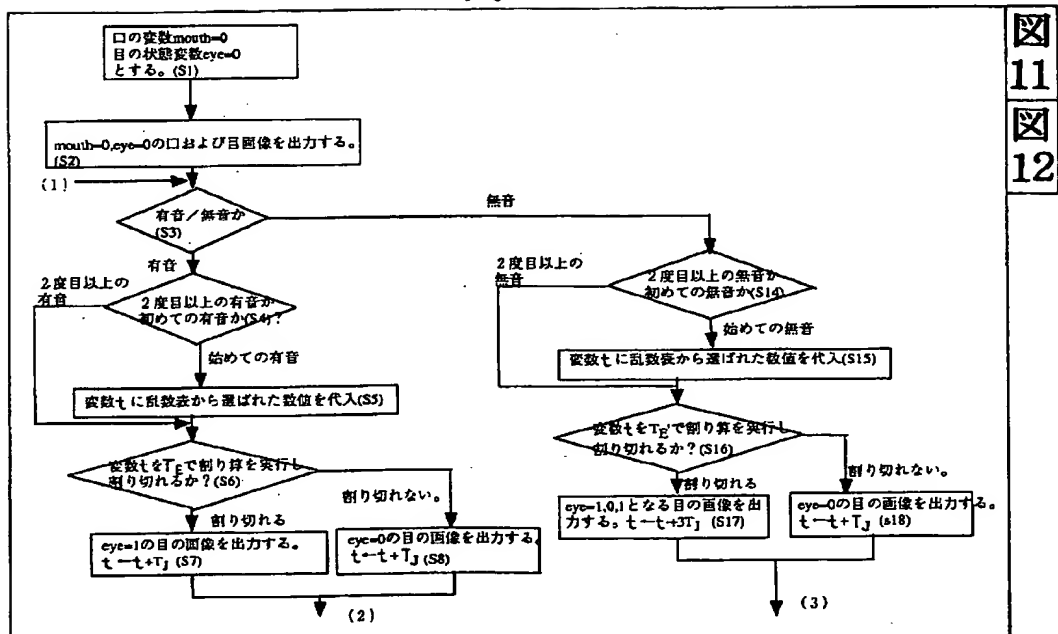
【図10】



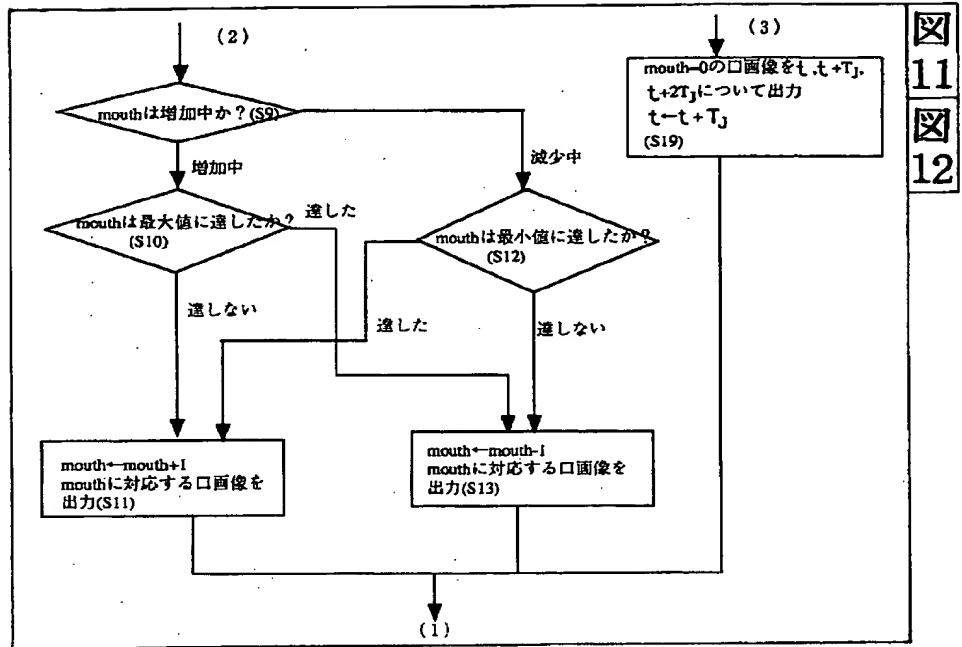
【図14】



【図11】



【図12】



フロントページの続き

(72)発明者 渡辺 信之  
東京都新宿区西新宿三丁目19番2号 日本  
電信電話株式会社内

(72)発明者 金山 英明  
東京都新宿区西新宿三丁目19番2号 日本  
電信電話株式会社内